



DUEL: A Multi-lingual Multimodal Dialogue Corpus for Disfluency, Exclamations and Laughter

Ye Tian, Julian Hough, Laura de Ruiter, Simon Betz, David Schlangen,
Jonathan Ginzburg

► To cite this version:

Ye Tian, Julian Hough, Laura de Ruiter, Simon Betz, David Schlangen, et al.. DUEL: A Multi-lingual Multimodal Dialogue Corpus for Disfluency, Exclamations and Laughter. LREC 2016, Tenth International Conference on Language Resources and Evaluation, 2016, Portorož,, Slovenia. hal-01371394

HAL Id: hal-01371394

<https://u-paris.hal.science/hal-01371394>

Submitted on 26 Sep 2016

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

DUEL: A Multi-lingual Multimodal Dialogue Corpus for Disfluency, Exclamations and Laughter

Ye Tian¹, Julian Hough², Laura de Ruiter³, Simon Betz²,
David Schlangen², Jonathan Ginzburg¹

¹Laboratoire de linguistique formelle, Université Paris-Diderot

²Dialogue Systems Group, Bielefeld University

³School of Psychological Sciences, University of Manchester

1 Introduction

Natural, spontaneous dialogue corpora are rich resources for a variety of linguistic research. In this paper, we present the DUEL (‘Disfluency, exclamations and laughter in dialogue’ (Ginzburg et al., 2014)) corpus, consisting of 24 hours of natural, face-to-face, loosely task-directed dialogue in German, French and Mandarin Chinese. The corpus is uniquely positioned as a cross-linguistic, multimodal dialogue resource controlled for domain, including audio, video and body tracking data and is transcribed and annotated for disfluency, laughter and exclamations.

To ensure cross-linguistic comparability, the experimental tasks were designed to be culture-neutral, the data in three languages were recorded using near-identical technical setups, and our transcription and annotation protocol is designed to be language-general.

In this paper, we give a summary of the tasks, the recording procedure and the transcription and annotation protocol. Then we discuss briefly the characteristics of our corpus and implications for natural dialogue research.

2 Existing Spontaneous Speech Corpora in Target Languages

Previous corpus work on spontaneous speech in German has focused small domains and/or on speech data that does not generalize well to natural face-to-face dialogue. Kohler (1996) elicited dialogues using an appointment making scenario, but had speakers press a button to speak, eliminating any turn-overlaps (and potential disfluencies resulting from these). This is similar to (Burger et al., 2000), who used a similar scenario and instructed speakers not to interrupt each other. Schiel et al. (2012)’s non-intoxicated spontaneous control data was obtained by having them talk to the experimenter in a car. Schmidt et al. (2010) and the

Berlin Map Task Corpus (BeMaTaC; <https://u.berlin.de/bematac>) both used map tasks, with the latter recording only non-native speakers. Peters (2005) collected a corpus of spontaneous speech by having two friends talk about video sequences via headset without eye-contact.

For French, there are several corpora for spontaneous speech. Several projects collected spoken French for studying prosody, for example, PFC (Durand et al., 2009), C-PROM (Avanzi et al., 2010) and Rhapsodie (Lacheret et al., 2014). Because of their research interests, these corpora cover a variety of discourse genres and do not focus on face-to-face dialogues. Bonneau-Maynard et al. (2005) collected the MEDIA corpus, containing roughly 70 hours of French dialogues in the topics of tourist information inquiry and hotel booking. It was recorded using a Wizard-of-Oz system where the participants interact with a human wizard they believe to be a machine. There are also corpora where the speech is not completely spontaneous, for example, the French oral narrative corpus (Carruthers, 2013) is a collection of stories told by storytellers.

For Mandarin Chinese, work on spontaneous speech is sparser. The NCCU corpus of spontaneous Chinese (Chui et al., 2008) contains face-to-face conversations (not necessarily between two speakers) in three languages: Mandarin, Hakka, and Southern Min. The Mandarin sub-corpus contains about 3.5 hours of conversations. The Lancaster Los Angeles Spoken Chinese Corpus (Xiao and Tao, 2007) is a collection of dialogues and monologues in Mandarin Chinese, both spontaneous and scripted. Recently, The Chinese Academy of Social Science initiated the on-going project “Spoken Chinese Corpus of Situated Discourse”, aiming to collect 1000 hours of spoken Chinese, covering different discourse genres and major dialects in China (cf.(Gu, 2000)).

The DUEL corpus is the first to provide French,

Chinese and German sub-corpora in comparable spontaneous dialogue domains with a unified disfluency and laughter mark-up, making it of potentially great interest to the dialogue and speech research communities.

3 Corpus Construction

We recorded 10 dyads per language. Each dyad participated in three tasks, with the whole interaction lasting roughly 45 minutes in total.

3.1 Task design

We devised the tasks with three goals in mind, for them to: 1) be specific enough so participants do not spend significant time in silence working out what they should do, but unconstrained enough to allow free speech; 2) help elicit laughter and exclamations (we assume that as long as the conversations are spontaneous, disfluencies occur regularly); and, 3) create different types of laughter depending on the nature of the roles the participants have in the tasks: laughter of pleasure (*Duchenne* laughter) and laughter of embarrassment and other interaction management (social laughter). The three tasks used were as follows:

Dream Apartment First used in (Kousidis et al., 2013), the participant pairs were told that they are to share a large open-plan apartment, and will receive a large amount of money (500,000 Euros) to furnish and decorate the apartment. The two participants are allowed their own bedroom but will share the rest of the apartment. They discuss the layout, furnishing and decoration of their apartment for 15 minutes.

Film Script This more open task requires the participants to spend 15 minutes creating a scene for a film in which something embarrassing happens to the main character. They are told that they can draw on their own experience.

Border Control This role-play interview task is the most constructed. One participant plays the role of a traveler attempting to pass through the border control of an imagined country, and is interviewed by an officer. The traveler has a personal history and situation that disfavors them in this interview (for example having a criminal record and carrying illegal substances). The officer asks questions that are general as well as specific to this traveler. In addition, the traveler happens to be

parent-in-law of the officer. For this task, the two participants receive separate information regarding their character roles, and the task is not timed – it ends when they feel that the interview is finished. The purpose of this task is to bring in an element of power asymmetry in the roles of the participants, while the other two can be considered symmetrical.

After the three tasks, the participants complete a questionnaire about the pair’s relationship (whether or for how long they know each other, and the frequency of contact) and how they felt about the tasks: how much they understood each other, and to what extent they felt uncomfortable or embarrassed during each task. This meta-data is available with the corpus.

3.2 Languages and participants

There were 10 pairs of native speakers for each of the three languages: German, French, and Chinese. The German speakers were all students at Bielefeld university where 3 pairs were friends/acquaintances and the remaining 7 strangers. The French speakers were students at Université Paris Diderot– 5 pairs were friends/acquaintances, and 5 were strangers. Among the 10 pairs of Chinese speakers, 7 pairs were university students in Paris, and 3 pairs were recruited via a local Chinese forum and, again, 5 pairs were friends/acquaintances and 5 were strangers. Participant gender was not controlled for.

3.3 Recording setup

The German data was recorded at Bielefeld University. The French and Chinese data were recorded at Université Paris Diderot.

The filming used two cameras to capture the gesture space and face of both participants, close lapel microphones were used to capture excellent audio quality without being intrusive, and the body movement was tracked by a Microsoft Kinect 2.

The body tracking data was logged into a time-stamped XML format using the *Venice.hub* logger (Kennington et al., 2014) which can be easily interpreted through the freely available Mumodo analysis tool kit.¹

¹Available from <https://github.com/dsg-bielefeld/mumodo>.

4 Transcription and Segmentation

Transcription was done from the WAV audio files using Praat (Boersma and Weenink, 2010), following the instructions of the DUEL manual (Hough et al., 2015). The manual specifies language general practices such as segmentation, disfluency annotation and laughter annotation, as well as language specific instructions regarding filled pauses, exclamations, and non-standard orthography.

For the transcription, the following tiers are available for a given participant X:

X-turns tier containing the turn boundaries for participant X

X-utts tier used for segmentation and transcription of X's utterances

X-en tier containing English paraphrase translation for X's utterances

Segmentation: In the tiers containing turns, all continuous stretches of speech by one speaker until the other speaker takes over, modulo small overlaps, are considered one turn. In the utterance tiers, we follow Meteer et al. (1995)'s notion of a *slash unit*, defining the notion of utterance as “maximally a sentence but can be a smaller unit [...] Intuitively, slash-units below the sentence level correspond to those parts of the narrative which are not sentential but which the annotator interprets as complete.”

5 Disfluency, Laughter and Exclamation Annotation

Our annotations follow the light-weight inline method of dialogue annotation described by Hough et al. (2015).

Disfluency: We consider disfluencies anything that leads to an audible deviation from expected speech production. We annotated the following phenomena: silent pauses, lengthening, filled pauses and editing terms, repairs, abandoned utterances and restarts.

For silent pauses, we transcribed pauses of short, medium and long duration, using one, two and three dots respectively. Lengthening was transcribed using the symbol “:” following the lengthened syllable(s), e.g. u:m:.

We mark filled pauses by a {F }, bracketing other fillers simply with { } - e.g. I { you know } like her.

The inventory of editing phrases and filled pauses differ depending on the language. For example, in German, the common filled pauses are {F äh}, {F ähm} and {F hm}; in French they are {F euh}, {F mmh} and {F euhm}; in Chinese, they are {F en}, {F eh}, as well as demonstratives {F nage} (literally “that”) and {F zhege} (literally “this”).

For repairs, restarts and abandoned utterances, we mark the structure according to this scheme (similar to (Meteer et al., 1995)):

(*reparandum* + { *editing term* } *repair*)

Both the editing term (which can be a filled pause) and the repair are optional. The structure can be nested and can appear in any positions in an utterance. Here are a few examples:

- (1) Standard repair: I went to (the: + {F um } the) garden
- (2) Nested: (I + (I + I)) want to go to Berlin
- (3) Restart: (I + {F uh }) yesterday someone said yes to that

For partial words, transcribers were encouraged to guess the complete standard form of the word where possible, again using a simple tag <p s=" . . "> . . -</p>, as below:

- (4) (<p s="Wohnzimmer">Wohn-</p> + . { ja also } (die + (die + das)) {F äh} ... Wohnzimmer)
(<p s="living room">liv-</p> { yes well } (the + (the + the)) {F uh } living room)

Laughter: We distinguish laughter concurrent with speech (laughed speech) and standalone laughter bouts. The former is transcribed with simple XML-style tags spanning the affected speech, e.g. <laughter>...</laughter>, and the latter is marked <laughter/>. A <laughterOffset/> tag as in (5) is used for the often audible deep inhalation of breath after laughed speech or a bout.

- (5) (Und mit einem +) mit vielleicht Sachen die nicht
<laughter> auseinander brechen </laughter>
<laughterOffset/> -
(And with a +) with perhaps things that don't
<laughter> fall apart </laughter>
<laughterOffset/> -

Example 1 (Chinese): Chaining repeat repair:

A	就	感觉	客厅	是	((((公+公)+公)+公共:))	{F jiushi }	休息	啊	或者
A-en	then	feel	living room	is	((((public+public)+public)+public:))	{F that is }	relax	PRT	or

Example 2 (German): Chaining substitution repair after laughed speech:

A	dann hat jeder genug Privatsphäre .. mit seinem <laughter> Partner </laughter>
	(und die Küche + (und die + {F ähm } (und die + ... und das Wohnzimmer))) ist quasi so ... mittig
A-en	then everyone has some privacy ... with their <laughter> partner </laughter>
	(and the kitchen + (and the + {F um } (and the + ... and the living room))) is kind of ... central

Example 3 (French): Restart disfluency within laughed speech:

A	bah quand même <laughter> c'est (un chien +) deuxième étage </laughter>
A-en	well still <laughter> it is (a dog+) second floor </laughter>

Figure 1: DUEL's disfluency and laughter mark-up in the three languages in the Dream Apartment task

Exclamations: We mark any exclamative short utterances with a simple XML-style tagging again, for example, <x>ohlala</x> in French. Compared to disfluencies and laughter, exclamations were sparse in our corpus, but the forms between languages is a fruitful area of cross-linguistic research.

Non-standard pronunciation: In the data of all three languages, there are pronunciations that deviate from the standard forms. Very often, the deviation is conventional. These are marked by transcribing both the actual form as well as the form in standard orthography as in <v s="standard form">pronounced form</v>. For example, <v s="auf der">aufer</v> (German), and <v s="il faut">'faut</v> (French). Similarly, for obvious mis-pronunciations we use <m s="standard form">pronounced form</m> (e.g. <m s="angry">angly</m>).

Non-verbal utterances For monosyllabic functional utterances we use hm (with no brackets) and for the disyllabic equivalent we use mhm. As regards other non-verbal vocalizations, apart from laughter and breathing and these functional utterances, for non-linguistic contributions such as coughing, sneezing and lip-smacking, we use a <nonverbal/> tag.

6 Discussion

DUEL's light-weight and consistent mark-up of the above phenomena allows for fast searching of the utterance tiers, and example utterances from the Dream Apartment task with the mark-up across the three languages can be seen in Figure 1. The mark-up exhibits good inter-annotator agreement and it is compatible with several existing schemes – see Hough et al. (2015) for details.

The audio and transcription and annotation data can be used in conjunction with the body-tracking data – the Dream Apartment task has been used in a study on the multimodal aspects of laughter by Kousidis et al. (2015) and continues to be used for multimodal dialogue studies. We will detail use cases in the full paper, should the paper be accepted.

7 Conclusion

We have presented the DUEL corpus, a multilingual, multimodal data-set that is uniquely positioned for dialogue and spontaneous speech research, both in terms of the consistency of the domain across languages, its standardization of disfluency and laughter mark-up and its synchronized multimodal data.

References

- Mathieu Avanzi, Anne-Catherine Simon, Jean-Philippe Goldman, and Antoine Auchlin. 2010. C-prom: An annotated corpus for french prominence study. In *Proceedings of Prosodic Prominence, Speech Prosody 2010 Workshop*.
- Paul Boersma and David Weenink. 2010. Praat: doing phonetics by computer.
- Hélène Bonneau-Maynard, Sophie Rosset, Christelle Ayache, Anne Kuhn, and Djamel Mostefa. 2005. Semantic annotation of the french media dialog corpus. In *Ninth European Conference on Speech Communication and Technology*.
- Susanne Burger, Karl Weilhammer, Florian Schiel, and Hans G Tillmann. 2000. Verbmobil data collection and annotation. In *Verbmobil: Foundations of Speech-to-Speech Translation*, pages 537–549. Springer.
- Janice Carruthers. 2013. French oral narrative corpus. Commissioning Body / Publisher: Oxford Text Archive.

- Kawai Chui, Huei-ling Lai, et al. 2008. The nccu corpus of spoken chinese: Mandarin, hakka, and southern min.
- Jacques Durand, Bernard Laks, and Chantal Lyche. 2009. Le projet pfc (phonologie du français contemporain): une source de données primaires structurées. *Phonologie, variation et accents du français*, pages 19–61.
- Jonathan Ginzburg, Ye Tian, Pascal Amsili, Claire Beyssade, Barbera Hemforth, Yannick Mathieu, Claire Saillard, Julian Hough, Spyros Kousidis, and David Schlangen. 2014. The Disfluency, Exclamation and Laughter in Dialogue (DUEL) Project. In *Proceedings of the 18th SemDial Workshop (DialWatt)*, pages 176–178, Herriot Watt University, Edinburgh.
- Yueguo Gu. 2000. Compiling a spoken chinese corpus of situated discourse. In *Keynote speech given at the 8th national conference on contemporary linguistics. Guangzhou*.
- Julian Hough, Laura de Ruiter, Simon Betz, and David Schlangen. 2015. Disfluency and laughter annotation in a light-weight dialogue mark-up protocol. In *The 6th Workshop on Disfluency in Spontaneous Speech (DiSS)*.
- Casey Kennington, Spyridon Kousidis, and David Schlangen. 2014. Multimodal dialogue systems with inprots and venice. In *Proceedings of the 18th SemDial Workshop on the Semantics and Pragmatics of Dialogue (DialWatt). Posters*.
- Klaus J Kohler. 1996. Labelled data bank of spoken standard german: the kiel corpus of read/spontaneous speech. In *ICSLP 96*, volume 3, pages 1938–1941. IEEE.
- Spyridon Kousidis, Thies Pfeiffer, and David Schlangen. 2013. Mint. tools: Tools and adaptors supporting acquisition, annotation and analysis of multimodal corpora. *Interspeech 2013*.
- Spyridon Kousidis, Julian Hough, and David Schlangen. 2015. Exploring the body and head kinematics of laughter, filled pauses and breaths. In *Proceedings of The 4th Interdisciplinary Workshop on Laughter and Other Non-verbal Vocalisations in Speech*, pages 23–25.
- Anne Lacheret, Sylvain Kahane, Julie Beliao, Anne Dister, Kim Gerdes, Jean-Philippe Goldman, Nicolas Obin, Paola Pietrandrea, Atanas Tchobanov, et al. 2014. Rhapsodie: a prosodic-syntactic treebank for spoken french. In *Language Resources and Evaluation Conference*.
- Marie W Meteer, Ann A Taylor, Robert MacIntyre, and Rukmini Iyer. 1995. *Disfluency annotation stylebook for the switchboard corpus*. University of Pennsylvania.
- B Peters. 2005. The database-the kiel corpus of spontaneous speech. *Prosodic Structures in German Spontaneous Speech, AIPUK 35a*, pages 1–6.
- Florian Schiel, Christian Heinrich, and Sabine Barfusser. 2012. Alcohol language corpus: the first public corpus of alcoholized german speech. *Language resources and evaluation*, 46(3):503–521.
- Thomas Schmidt, Hanna Hedeland, Timm Lehmborg, and Kai Wörner. 2010. Hamatac—the hamburg map-task corpus.
- Richard Xiao and Hongyin Tao. 2007. The lancaster los angeles spoken chinese corpus.